## BRIEF DESCRIPTION OF THE DRAWINGS

[0009] FIG. 1 is a cartoon illustration of an exemplary vector of the present disclosure useful for generating a double-stranded nucleic acid library.

[0010] FIG. 2 is a cartoon illustration of an exemplary vector of the present disclosure, wherein adaptor sequences are included and are useful for, for example, bridge amplification methods before sequencing.

[0011] FIGS. 3A and 3B show characteristics of a cypher library and the detection of true mutations. (A) Data generated in a single next generation sequence run on MiSeq® demonstrates broad coverage and diversity at the upstream seven base pair cypher in a vector library, wherein the vector used is illustrated in FIG. 2. (B) Cypher Seq eliminates errors introduced during library preparation and sequencing. Target nucleic acid molecules were ligated into a cypher vector library containing previously catalogued dual, double-stranded cyphers. The target sequences were amplified and sequenced. All sequencing reads having identical cypher pairs, along with their reverse complements, were grouped into families. Comparison of family sequences allowed for generation of a consensus sequence wherein 'mutations' (errors) arising during library preparation (open circle) and during sequencing (gray circle and triangle) were computationally eliminated. Generally, mutations that are present in all or substantially all reads (black diamond) from the same cypher and its reverse complement are counted as true mutations.

[0012] FIGS. 4A and 4B show that the cypher system can distinguish true mutations from artifact mutations. (A) Wild-type TP53 Exon 4 was ligated into a library of Cypher Seq vectors and sequenced on the Illumina MiSeq® instrument with a depth of over a million. Sequences were then compared to wild-type TP53 sequence. Detected substitutions were plotted before (A) and after correction (B) with Cypher Seq.

## DETAILED DESCRIPTION

[0013] In one aspect, the present disclosure provides a double-stranded nucleic acid library wherein target nucleic acid molecules include dual cyphers (i.e., barcodes or origin identifier tags), one on each end (same or different), so that sequencing each complementary strand can be connected or linked back to the original molecule. The unique cypher on each strand links each strand with its original complementary strand (e.g., before any amplification), so that each paired sequence serves as its own internal control. In other words, by uniquely tagging double-stranded nucleic acid molecules, sequence data obtained from one strand of a single nucleic acid molecule can be specifically linked to sequence data obtained from the complementary strand of that same double-stranded nucleic acid molecule. Furthermore, sequence data obtained from one end of a double-stranded target nucleic acid molecule can be specifically linked to sequence data obtained from the opposite end of that same double-stranded target nucleic acid molecule (if, for example, it is not possible to obtain sequence data across the entire target nucleic acid molecule fragment of the library).

[0014] The compositions and methods of this disclosure allow a person of ordinary skill in the art to more accurately distinguish true mutations (i.e., naturally arising in vivo mutations) of a nucleic acid molecule from artifact "mutations" (i.e., ex vivo mutations or errors) of a nucleic acid molecule that may arise for various reasons, such as a downstream amplification error, a sequencing error, or physical or chemical damage. For example, if a mutation pre-existed in the original double-stranded nucleic acid molecule before isolation, amplification or sequencing, then a transition mutation of adenine (A) to guanine (G) identified on one strand will be complemented with a thymine (T) to cysteine (C) transition on the other strand. In contrast, artifact "mutations" that arise later on an individual (separate) DNA strand due to polymerase errors during isolation, amplification or sequencing are extremely unlikely to have a matched base change in the complementary strand. The approach of this disclosure provides compositions and methods for distinguishing systematic errors (e.g., polymerase read fidelity errors) and biological errors (e.g., chemical or other damage) from actual known or newly identified true mutations or single nucleotide polymorphisms (SNPs).

[0015] In certain embodiments, the two cyphers on each target molecule have sequences that are distinct from each other and, therefore, provide a unique pair of identifiers wherein one cypher identifies (or is associated with) a first end of a target nucleic acid molecule and the second cypher identifies (or is associated with) the other end of the target nucleic acid molecule. In certain other embodiments, the two cyphers on each target molecule have the same sequence and, therefore, provide a unique identifier for each strand of the target nucleic acid molecule. Each strand of the double-stranded nucleic acid library (e.g., genomic DNA, cDNA) can be amplified and sequenced using, for example, next generation sequencing technologies (such as, emulsion PCR or bridge amplification combined with pyrosequencing or sequencing by synthesis, or the like). The sequence information from each complementary strand of a first double-stranded nucleic acid molecule can be linked and compared (e.g., computationally "de-convoluted") due to the unique cyphers associated with each end or strand of that particular double-stranded nucleic acid molecule. In other words, each original double-stranded nucleic acid molecule fragment found in a library of molecules can be individually reconstructed due to the presence of an associated unique barcode or pair of barcode (identifier tag) sequences on each target fragment or strand.

[0016] By way of background, any spontaneous or induced mutation will be present in both strands of a native genomic, double-stranded DNA molecule. Hence, such a mutant DNA template amplified using PCR will result in a PCR product in which 100% of the molecules produced by PCR include the mutation. In contrast to an original, spontaneous mutation, a change due to polymerase error will only appear in one strand of the initial template DNA molecule (while the other strand will not have the artifact mutation). If all DNA strands in a PCR reaction are copied equally efficiently, then any polymerase error that emerges from the first PCR cycle likely will be found in at least 25% of the total PCR product. But DNA molecules or strands are not copied equally efficiently, so DNA sequences amplified from the strand that incorporated an erroneous nucleotide base during the initial amplification might constitute more or less than 25% of the population of amplified DNA sequences depending on the efficiency of amplification, but still far less than 100%. Similarly, any polymerase error that occurs in later PCR cycles will generally represent an even smaller proportion of PCR products (i.e., 12.5% for the second